# Similarity Search The q-Gram Distance

#### Nikolaus Augsten

nikolaus.augsten@plus.ac.at Department of Computer Science University of Salzburg



WS 2025/26

Version November 6, 2025

- Filters for the Edit Distance
  - Motivation
  - Lower Bound Filters
  - Length Filter
  - q-Grams: Count Filter
  - q-Grams: Position Filtering
  - Experiments
- The q-Gram Distance
- 3 Conclusion

- Filters for the Edit Distance
  - Motivation
  - Lower Bound Filters
  - Length Filter
  - q-Grams: Count Filter
  - q-Grams: Position Filtering
  - Experiments
- 2 The q-Gram Distance
- 3 Conclusion

- Filters for the Edit Distance
  - Motivation
  - Lower Bound Filters
  - Length Filter
  - q-Grams: Count Filter
  - q-Grams: Position Filtering
  - Experiments
- 2 The q-Gram Distance
- 3 Conclusion

## Application Scenario

- Scenario:
  - A company offers a number of services on the Web.
  - You can subscribe for each service independently.
  - Each service has its own database (no unique key across databases).
- Example: customer tables of two different services:

Α				В	
ID	name		ID	name	
1023	Frodo Baggins		948483	John R. R. Tolkien	
21	J. R. R. Tolkien		153494	C. S. Lewis	
239	C.S. Lewis		494392	Fordo Baggins	
863	Bilbo Baggins		799294	Biblo Baggins	

Task: Created unified customer view!

## The Join Approach

• Solution: Join customer tables on name attribute (Q1):

```
SELECT * FROM A,B
WHERE A.name = B.name
```

- Exact Join: Does not work!
- Similarity Join: Allow *k* errors...
  - (1) Register UDF (User Defined Function) for the edit distance:

returns the union cost edit distance between the strings x and y.

(2) Rewrite query Q1 as similarity join (Q2):

```
SELECT * FROM A,B
WHERE ed(A.name, B.name) <= k</pre>
```

# Effectiveness and Efficiency of the Approximate Join

• Effectiveness: Join result for k = 3:

ID	name	ID	name
1023	Frodo Baggins	494392	Fordo Baggins
21	J. R. R. Tolkien	948483	John R. R. Tolkien
239	C.S. Lewis	153494	C. S. Lewis
863	Bilbo Baggins	799294	Biblo Baggins

- $\Rightarrow$  very good (100% correct)
- Efficiency: How does the DB evaluate the query?
  - (1) compute  $A \times B$
  - (2) evaluate UDF on each tuple  $t \in A \times B$
- Prohibitive runtime!

# Using a Filter for Search Space Reduction

- Search space:  $A \times B \ (\Rightarrow |A| \cdot |B| \ \text{edit distance computations})$
- Filtering (Pruning): Remove tuples that can not match, without actually computing the distance.

- Filters for the Edit Distance
  - Motivation
  - Lower Bound Filters
  - Length Filter
  - q-Grams: Count Filter
  - q-Grams: Position Filtering
  - Experiments
- 2 The q-Gram Distance
- 3 Conclusion

Carret Desult

## Filter Properties

• Error Types:

		Correct Result			
		positive	negative		
Filter	positive	true positive	false positive		
Test	negative	false negative	true negative		

- Example: "Are x and y within edit distance k?"
  - Correct result: compute edit distance and test  $ed(x, y) \le k$
  - Filter test: give answer without computing edit distance
  - False negatives: x and y are pruned although  $ed(x, y) \le k$ .
  - False positives: x and y are not pruned although  $ed(x, y) \nleq k$ .
- Good filters have
  - no false negatives (i.e., miss no correct results)
  - few false positive (i.e., avoid unnecessary distance computations)

#### Lower Bound Filters

• Lower bound (lb) for distance dist(x, y):

$$dist(x, y) \ge Ib_{dist}(x, y)$$

Query Q3 with Lower Bound Filter :

```
SELECT * FROM A,B
WHERE lb(A.name, B.name) <= k AND
ed(A.name, B.name) <= k</pre>
```

- lb(A.name, B.name) is a cheap function
- database will optimize query: compute ed(A.name,B.name) only if lb(A.name,B.name)  $\leq k$
- No false negatives!

- Filters for the Edit Distance
  - Motivation
  - Lower Bound Filters
  - Length Filter
  - q-Grams: Count Filter
  - q-Grams: Position Filtering
  - Experiments
- 2 The q-Gram Distance
- 3 Conclusion

## Length Filtering

#### Theorem (Length Filtering $[GIJ^+01]$ )

If two strings x and y are within edit distance k, their lengths cannot differ by more than k:

$$\operatorname{ed}(x,y) \geq \operatorname{abs}(|x| - |y|)$$

- Proof: At least abs(|x| |y|) inserts are needed to bring x and y to the same length.
- Query Q4 with Length Filtering:

```
SELECT * FROM A,B
WHERE ABS(LENGTH(A.name)-LENGTH(B.name)) <= k AND
ed(A.name, B.name) <= k</pre>
```

# Example: Length Filtering

• Execute query without/with length filter (k = 3):

A			В		
ID	name		ID	name	
1023	Frodo Baggins <sub>13</sub>	•	948483	John R. R. Tolkien <sub>18</sub>	
21	J. R. R. Tolkien <sub>16</sub>		153494	C. S. Lewis <sub>11</sub>	
239	C.S. Lewis <sub>10</sub>		494392	Fordo Baggins <sub>13</sub>	
863	Bilbo Baggins <sub>13</sub>		799294	Biblo Baggins <sub>13</sub>	

- Without length filter: 16 edit distance computations
- With length filter (k = 3): 12 edit distance computations
  - ullet J. R. R. Tolkien  $\leftrightarrow$  C. S. Lewis is pruned
  - all pairs (..., John R. R. Tolkien) except (J. R. R. Tolkien, John R. R. Tolkien) are pruned

- Filters for the Edit Distance
  - Motivation
  - Lower Bound Filters
  - Length Filter
  - q-Grams: Count Filter
  - q-Grams: Position Filtering
  - Experiments
- 2 The q-Gram Distance
- 3 Conclusion

# What is a q-Gram?

- Intuition:
  - slide window of length q over string  $x \in \Sigma^*$
  - characters covered by window form a q-gram
  - where window extends string: fill with dummy character #  $\notin \Sigma$

```
• Example: x = \text{Frodo}, q = 3
  extended: ##Frodo##
  q-grams: ##F
             #Fr
               Fro
                rod
                  o d o
                   d o #
                     0##
```

- q-Gram Profile  $G_x$ : bag of all q-grams of x
- Profile size:  $|G_x| = |x| + q 1$

# Single Edit Operations and Changing q-Grams

- Intuition: Strings within small edit distance share many q-grams.
- How many q-grams (q = 3) change/remain?

X	$ G_x $	У	$ G_y $	$ G_x \cap G_y $
peter	7	meter	7	4
peter	7	peters	8	5
peter	7	peer	6	4

 $\bullet$  ed $(x, y) = 1 \Rightarrow |G_x \cap G_y| = \max(|G_x|, |G_y|) - q$ 

# Multiple Edit Operations and Changing q-Grams

- $\bullet$  ed $(x,y) = 1 \Rightarrow |G_x \cap G_y| = \max(|G_x|,|G_y|) q$
- What if ed(x, y) = k > 1?

X	$ G_x $	y	$ G_y $	$ G_x \cap G_y $
peter	7	meters	8	2
peter	7	petal	7	3

• Multiple edit operations may affect the same q-gram:

$$peter \rightarrow G_x = \{ \#p, \#pe, pet, ete, ter, er\#, r\#\# \}$$
 $petal \rightarrow G_y = \{ \#p, \#pe, pet, eta, tal, al\#, l\#\# \}$ 

Each edit operation affects at most q q-grams.

## Count Filtering

#### Theorem (Count Filtering [GIJ $^+$ 01])

Consider two strings x and y with the q-gram profiles  $G_x$  and  $G_y$ , respectively. If x and y are within edit distance k, then the cardinality of the q-gram profile intersection is at least

$$|G_X \cap G_V| \ge \max(|G_X|, |G_V|) - kq$$

- Proof (by induction):
  - true for k=1:  $|G_x \cap G_y| \ge \max(|G_x|, |G_y|) q$
  - $k \to k+1$ : each additional edit operation changes at most q q-grams.  $\square$

## Implementation of q-Grams

- Given: tables A and B with schema (id, name)
  - *id* is the key attribute
  - name is string-valued
- Compute auxiliary tables QA and QB with schema (id, qgram):
  - each tuple stores one q-gram
  - string x of attribute name is represented by its |x| + q 1 q-grams
  - QA.id is the key value (A.id) of a tuple with A.name = x
  - QA.qgram is one of the q-grams of x
- Example:

	$\mathcal{A}$	(	QA
id	name	id	qgram
1023	Frodo Baggins	1023	##F
21	J. R. R. Tolkien	1023	#Fr
239	C.S. Lewis		
863	Bilbo Baggins	21	##J
	30	21	<b>#</b> J.

# Count Filtering Query

• Query Q5 with Count Filtering:

```
SELECT A.id, B.id, A.name, B.name

FROM A, QA, B, QB

WHERE A.id = QA.id AND

B.id = QB.id AND

QA.qgram = QB.qgram AND

ABS(LENGTH(A.name)-LENGTH(B.name)) <= k

GROUP BY A.id, B.id, A.name, B.name

HAVING COUNT(*) >= LENGTH(A.name)-1-(k-1)*q AND

COUNT(*) >= LENGTH(B.name)-1-(k-1)*q AND

ed(A.name,B.name) <= k
```

# Problem with Count Filtering Query

- Previous query Q5 works fine for  $kq < \max(|G_x|, |G_y|)$ .
- However: If  $kq \ge \max(|G_x|, |G_y|)$ , no q-grams may match even if  $\operatorname{ed}(x, y) <= k$ .
- False negatives:
  - short strings with respect to edit distance (e.g., |x| = 3, k = 3)
  - even if within given edit distance, matches tend to be meaningless (e.g., abc and xyz are within edit distance k = 3)

## Fixing Count Filtering Query

- Fix query to avoid false negatives [GIJ<sup>+</sup>03]:
  - Join pairs (x, y) with  $kq \ge \max(|G_x|, |G_y|)$  using only length filter.
  - Union results with results of previous query Q5.
- Query Q6 without false negatives (extends previous query Q5):

```
UNION
SELECT A.id, B.id, A.name, B.name
FROM
     А, В
WHERE LENGTH(A.name)+q-1 \le k*q AND
       LENGTH(B.name)+q-1 \le k*q AND
       ABS(LENGTH(A.name) - LENGTH(B.name)) <= k AND
       ed(A.name, B.name) <= k
```

 Note: We omit this part in subsequent versions of the query since it remains unchanged.

- Filters for the Edit Distance
  - Motivation
  - Lower Bound Filters
  - Length Filter
  - q-Grams: Count Filter
  - *q*-Grams: Position Filtering
  - Experiments
- 2 The q-Gram Distance
- 3 Conclusion

## Positional q-Grams

- Enrich *q*-grams with position information:
  - extended string: prefix and suffix string x with q-1 characters #
  - slide window of length q over extended string x'
  - characters covered by window after shifting it i times form the q-gram at position i+1
- Example: x = Frodo##Frodo## extended string: positional q-grams: (1, # F)(2, # Fr)(3,Fro)(4,r od)(5, odo)(6,do#)(7,0 ##)

## Computing Positional q-Grams in SQL

- Given: table N
  - N has a single attribute i
  - N is filled with numbers from 1 to max (max is the maximum string length plus q-1)
- Positional q-grams for table A in SQL (Q7):

```
CREATE TABLE QA AS
       SELECT A.id, N.i AS pos,
           SUBSTRING(CONCAT(
               SUBSTRING('\#..\#', 1, q - 1),
               LOWER(A.name),
               SUBSTRING('\#..\#', 1, q - 1)),
           N.i, q) AS qgram
       FROM A, N
       WHERE N.i <= LENGTH(A.name) + q - 1
```

## Corresponding *q*-Grams

- Corresponding *q*-gram:
  - Given: positional q-grams (i, g) of x
  - transform x to y applying edit operations
  - (i,g) "becomes" (j,g) in y
  - We define: (i,g) corresponds to (j,g)
- Example:
  - x' = #abaZabaabaaba##, y' = #abaabaabaabaabaaba##
  - edit distance is 1 (delete Z from x)
  - (7, aba) in x corresponds to (6, aba) in y
  - ... but not to (9, <u>aba</u>)

# Position Filtering

#### Theorem (Position Filtering [GIJ<sup>+</sup>01])

If two strings x and y are within edit distance k, then a positional q-gram in one cannot correspond to a positional q-gram in the other that differs from it by more then k positions.

- Proof:
  - each increment (decrement) of a position requires an insert (delete);
  - a shift by k positions requires k inserts/deletes.

## Position Filtering

Query Q8 with Count and Position Filtering:

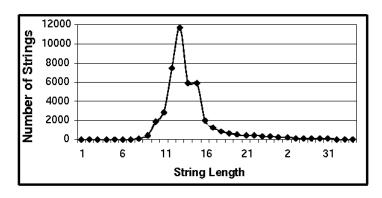
```
SELECT A.id, B.id, A.name, B.name
FROM A, QA, B, QB
WHERE A.id = QA.id AND
        B.id = QB.id AND
         QA.qgram = QB.qgram AND
         ABS(LENGTH(A.name)-LENGTH(B.name)) <= k AND
         ABS(QA.pos-QB.pos) <= k
GROUP BY A.id, B.id, A.name, B.name
        COUNT(*) >= LENGTH(A.name)-1-(k-1)*q AND
HAVING
         COUNT(*) >= LENGTH(B.name)-1-(k-1)*q AND
         ed(A.name, B.name) <= k
```

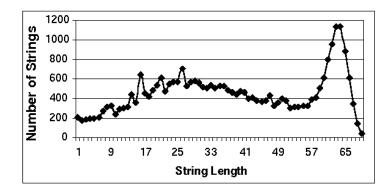
- Filters for the Edit Distance
  - Motivation
  - Lower Bound Filters
  - Length Filter
  - q-Grams: Count Filter
  - q-Grams: Position Filtering
  - Experiments
- 2 The q-Gram Distance
- 3 Conclusion

# Experimental Data

- All experimental results taken from [GIJ+01]
- Three string data sets:
  - set1: 40K tuples, average length: 14 chars
  - set2: 30K tuples, average length: 38 chars
  - set3: 30K tuples, average length: 33 chars

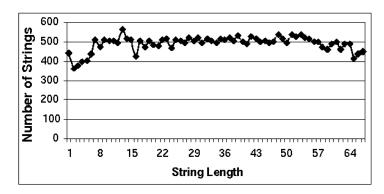
# String Length Distributions





Set 1

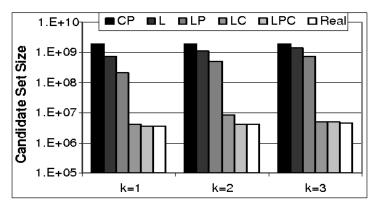
Set 2



Set 3

#### Candidate Set Size

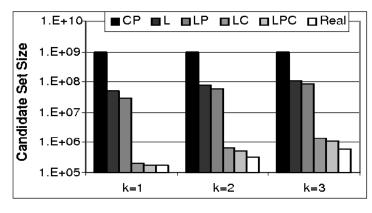
- Question: How many edit distances do we have to compute?
- Show candidate set size for different filters (small is good).
- q = 2
- Caption:
  - CP: cross product
  - L: length filtering, P: position filtering, C: count filtering
  - Real: number of real matches



Set 1

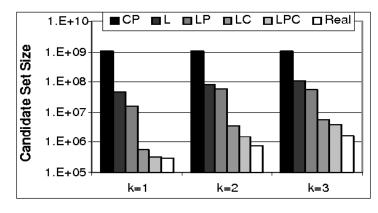
#### Candidate Set Size

- Question: How many edit distances do we have to compute?
- Show candidate set size for different filters (small is good).
- q = 2
- Caption:
  - CP: cross product
  - L: length filtering, P: position filtering, C: count filtering
  - Real: number of real matches



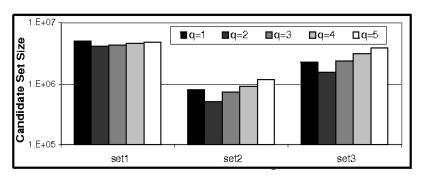
#### Candidate Set Size

- Question: How many edit distances do we have to compute?
- Show candidate set size for different filters (small is good).
- q = 2
- Caption:
  - CP: cross product
  - L: length filtering, P: position filtering, C: count filtering
  - Real: number of real matches

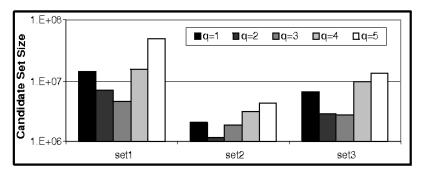


# Various q-Gram Lengths

- Question: How does the choice of q influence the filter effectiveness?
- Show candidate set size for different q values (small is good).



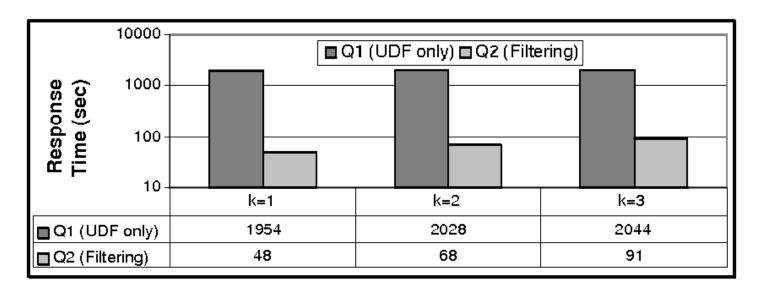
Edit Distance Threshold k = 2



Edit Distance Threshold k = 3

## Response Time

- Approximate self-join on small sample of 1000 tuples (set 1) (full dataset > 3 days without filters!)
- Measure response time (small is good).
- Caption:
  - k: edit distance threshold
  - Q1: edit distance without filters
  - Q2: edit distance with filters



#### Outline

- Filters for the Edit Distance
  - Motivation
  - Lower Bound Filters
  - Length Filter
  - q-Grams: Count Filter
  - q-Grams: Position Filtering
  - Experiments
- The q-Gram Distance
- 3 Conclusion

## The q-Gram Distance

#### Definition (q-Gram Distance [Ukk92])

Let  $G_x$  and  $G_y$  be the q-gram profiles of the strings x and y, respectively. The q-gram distance between two strings is the number of q-grams in  $G_x$  and  $G_y$  that have no match in the other profile,

$$\operatorname{dist}_{q}(x,y) = |G_{x} \uplus G_{y}| - 2|G_{x} \cap G_{y}|.$$

• Example: q = 2, x = abab, y = abcab

$$G_X = \{ \#a, ab, ba, ab, b\# \}$$
 $G_Y = \{ \#a, ab, bc, ca, ab, b\# \}$ 
 $G_X \uplus G_Y = \{ \#a, ab, ba, ab, b\#, \#a, ab, bc, ca, ab, b\# \}$ 
 $G_X \cap G_Y = \{ \#a, ab, ab, b\# \}$ 
 $dist_g(x, y) = |G_X \uplus G_Y| - 2|G_X \cap G_Y| = 11 - 2 \cdot 4 = 3$ 

## Pseudo Metric q-Gram Distance

- The q-gram distance is a pseudo metric: For all  $x, y, z \in \Sigma^*$ 
  - $\operatorname{dist}_q(x,y) + \operatorname{dist}_q(y,z) \ge \operatorname{dist}_q(x,z)$  (triangle inequality)
  - $\operatorname{dist}_q(x,y) = \operatorname{dist}_q(y,x)$  (symmetric)
  - $\operatorname{dist}_{q}(x, y) = 0 \Leftarrow x = y$
- Note: Identity condition relaxed:  $dist_q(x, y) = 0 \Rightarrow x = y$  i.e., the *q*-gram distance between two *different* strings can be 0
- Example:

$$dist_q(axybxycxyd, axycxybxyd) = 0$$

$$G_x = G_y = \{\#\#a, \#ax, axy, xyb, ybx, bxy, xyc, ycx, cxy, xyd, yd\#, d\#\#\}$$

# Distance Normalization (1/3)

• What is a good threshold?

- Problem: Absolute numbers not always meaningful...
- Solution: Compute error relative to string length!

# Distance Normalization (2/3)

- Normalize distance such that  $\delta(x, y) \in [0..1]$
- Edit Distance:
  - non-normalized:  $0 \le \operatorname{ed}(x, y) \le \max(|x|, |y|)$
  - normalized edit distance:  $0 \le \text{norm-ed}(x, y) \le 1$

$$norm-ed(x,y) = \frac{ed(x,y)}{max(|x|,|y|)}$$

- *q*-Gram Distance:
  - non-normalized:  $0 \le \operatorname{dist}_q(x,y) \le |G_x \uplus G_y| |G_x \cap G_y|$
  - normalized q-gram distance:  $0 \le \text{norm-dist}_q(x, y) \le 1$

$$\mathsf{norm\text{-}dist}_q(x,y) = \frac{\mathsf{dist}_q(x,y)}{|G_x \uplus G_y| - |G_x \cap G_y|} = 1 - \frac{|G_x \cap G_y|}{|G_x \uplus G_y| - |G_x \cap G_y|}$$

<sup>&</sup>lt;sup>1</sup>Jaccard normalization. Dividing by  $|G_x \uplus G_y|$  (Dice normalization) also normalizes to [0..1], but the metric properties (triangle inequality) get lost [ABG10].

# Distance Normalization (3/3)

Normalized edit distance:

```
\label{eq:norm-ed} \begin{array}{ll} \text{norm-ed(International Business Machines Corporation}, \\ & \text{International Bussiness Machine} \ \text{Corporation}) = 0.047 \\ \text{norm-ed(IBM, BMW)} = 0.66 \\ \text{norm-ed(Int. Business Machines Corp.}, \\ & \text{International Business Machines Corporation}) = 0.4 \\ \end{array}
```

• Normalized q-gram distance (q = 3):

```
\begin{array}{c} \text{norm-dist}_q(\text{International Business Machines Corporation},\\ \text{International Bussiness Machine Corporation}) = 0.089\\ \text{norm-dist}_q(\text{IBM},\text{BMW}) = 1.0\\ \text{norm-dist}_q(\text{Int. Business Machines Corp.},\\ \text{International Business Machines Corporation}) = 0.36 \end{array}
```

## Edit Distance vs. q-Gram Distance

• Edit distance can not handle block-moves well:

```
x = 	ext{Nikolaus Augsten} y = 	ext{Augsten Nikolaus} norm-ed(x,y) = 1.0 norm-dist_q(x,y) = 0.39 (q = 3)
```

• *q*-Gram distance may be too strict:

```
x = +39-06-46-74-22 y = (39\ 06\ 467422)
norm-ed(x, y) = 0.4
norm-dist_q(x, y) = 1.0 (q = 3)
```

#### Outline

- Filters for the Edit Distance
  - Motivation
  - Lower Bound Filters
  - Length Filter
  - q-Grams: Count Filter
  - q-Grams: Position Filtering
  - Experiments
- 2 The q-Gram Distance
- 3 Conclusion

## Summary

- Approximate join with edit distance inefficient.
- Edit distance filters speed up join:
  - Length filter: based on the string length
  - Count filter: based on *q*-Grams
  - Position filter: based on positional q-Grams

- Nikolaus Augsten, Michael Böhlen, and Johann Gamper.
  The pq-gram distance between ordered labeled trees.

  ACM Transactions on Database Systems (TODS), 35(1):1–36, 2010.
- Luis Gravano, Panagiotis G. Ipeirotis, H. V. Jagadish, Nick Koudas, S. Muthukrishnan, and Divesh Srivastava.

  Approximate string joins in a database (almost) for free.

  In *Proceedings of the International Conference on Very Large Databases (VLDB)*, pages 491–500, Roma, Italy, September 2001.

  Morgan Kaufmann Publishers Inc.
- Luis Gravano, Panagiotis G. Ipeirotis, H. V. Jagadish, Nick Koudas, S. Muthukrishnan, and Divesh Srivastava.

  Approximate string joins in a database (almost) for free Erratum. Technical Report CUCS-011-03, Department of Computer Science, Columbia University, 2003.
- Esko Ukkonen.

  Approximate string-matching with *q*-grams and maximal matches.

